

Varying Receptive Fields for Thermal Image Super-Resolution

Nolan B. Gutierrez and William J. Beksi

Department of Computer Science and Engineering, University of Texas at Arlington

The Goal

Determine a mapping between a low-resolution (LR) image and its high-resolution (HR) counterpart in the HR space.

Introduction

Advances in thermographic cameras have spurred research in the thermal image domain for a variety of applications, but the inherent LR of thermal imaging systems necessitates image restoration techniques. Super-resolution (SR) is an image restoration method which involves the discovery of a mapping from low-quality to high-quality images in the infrared image space. Besides applications in the visible spectrum, there are many uses of SR in the infrared (invisible) range including building inspection, military operations, maritime search and rescue, and much more.

Thermal cameras are subject to deteriorating computing conditions as the degradation of LR images increases thus requiring computationally efficient SR techniques. However, SR algorithms such as convolutional enhancements and visual attention often suffer from performance losses. Dilated convolutions can augment deep SR techniques by increasing receptive fields at a non-increasing computational complexity. By parametrizing several convolutional filters' dilation rate, we can sample from feature maps using distinct receptive fields. Furthermore, we can apply second order channel-attention to efficiently and informatively attend to easily-computed global statistics (e.g., covariance matrices) of feature channels.

Background

Assume $F(\mathbf{x}; \boldsymbol{\theta}) : \mathbf{x} \in \mathcal{R}^{H \times W \times c} \rightarrow \mathbf{y} \in \mathcal{R}^{sH \times sW \times c}$, where H , W , and c are the height, weight, and number of channels of an input image. \mathbf{x} , \mathbf{y} , $\boldsymbol{\theta}$, and s are the input image, super-resolved image, parameters associated with F , and scale factor marking the increase of resolution, respectively. We use the following objective (mean squared error)

$$\boldsymbol{\theta} = \arg \min_{\boldsymbol{\theta}} \frac{1}{NHW} \sum_{\mathbf{x} \in \mathbf{X}} \sum_{i=1}^N (F(\mathbf{x}; \boldsymbol{\theta})(i) - \mathbf{y}(i))^2 \quad (1)$$

to find $\boldsymbol{\theta}$, where N is the number of sample input images in a batch \mathbf{X} . $F(\mathbf{x}; \boldsymbol{\theta})(i)$ and $\mathbf{y}(i)$ are the i th pixels of the super-resolved and ground-truth images.

Proposed Method

Dilated Convolutions: As illustrated in the different dilation rates module of Fig. 1, we parametrize parallel convolutional filters thereby introducing distinct dilation rates.

Compression Through Dilations:

An *effective receptive field* (ERF) is defined as the area containing any input pixel with at least some impact on a particular output unit.

Compression through dilations (CTD) is the case in which a convolutional filter uses fewer parameters to increase an ERF with dilated convolutions compared to without dilated convolutions.

Second Order Channel-Attention (SOCA): First, a feature map of dimension $H \times W \times C$ is reshaped into a feature map \mathbf{X} of shape $HW \times C$. Second, the covariance matrix is calculated,

$$\boldsymbol{\Sigma} = \mathbf{X} \mathbf{I}_f \mathbf{X}^T, \quad (2)$$

where $\mathbf{I}_f = \frac{1}{s}(\mathbf{I} - \frac{1}{s}\mathbf{1})$ and $s = HW$. \mathbf{I} and $\mathbf{1}$ are the $m \times m$ identity matrix and the matrix of all ones, respectively. Next, the covariance matrix is pre-normalized,

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{\text{tr}(\boldsymbol{\Sigma})} \boldsymbol{\Sigma}, \quad (3)$$

where $\text{tr}(\cdot)$ denotes the matrix trace. Let $\mathbf{Y}_0 = \hat{\boldsymbol{\Sigma}}$ and $\mathbf{Z}_0 = \mathbf{I}$, then \mathbf{Y}_n and \mathbf{Z}_n are obtained by

$$\mathbf{Y}_n = \frac{1}{2} \mathbf{Y}_{n-1} (3\mathbf{I} - \mathbf{Z}_{n-1} \mathbf{Y}_{n-1}), \quad (4)$$

$$\mathbf{Z}_n = \frac{1}{2} (3\mathbf{I} - \mathbf{Z}_{n-1} \mathbf{Y}_{n-1}) \mathbf{Z}_{n-1}, \quad (5)$$

with \mathbf{Y}_n and \mathbf{Z}_n quadratically converging to \mathbf{Y} and \mathbf{Y}^{-1} , respectively.

The final normalized matrix after five iterations of Newton-Schulz is found by compensating the pre-normalization step with

$$\hat{\mathbf{Y}} = \sqrt{\text{tr}(\boldsymbol{\Sigma})} \mathbf{Y}_N. \quad (6)$$

Afterwards, global covariance pooling is applied to obtain a scalar-valued statistic z_i for each channel i ,

$$z_i = \frac{1}{C} \sum_j \hat{\mathbf{Y}}_{ij}. \quad (7)$$

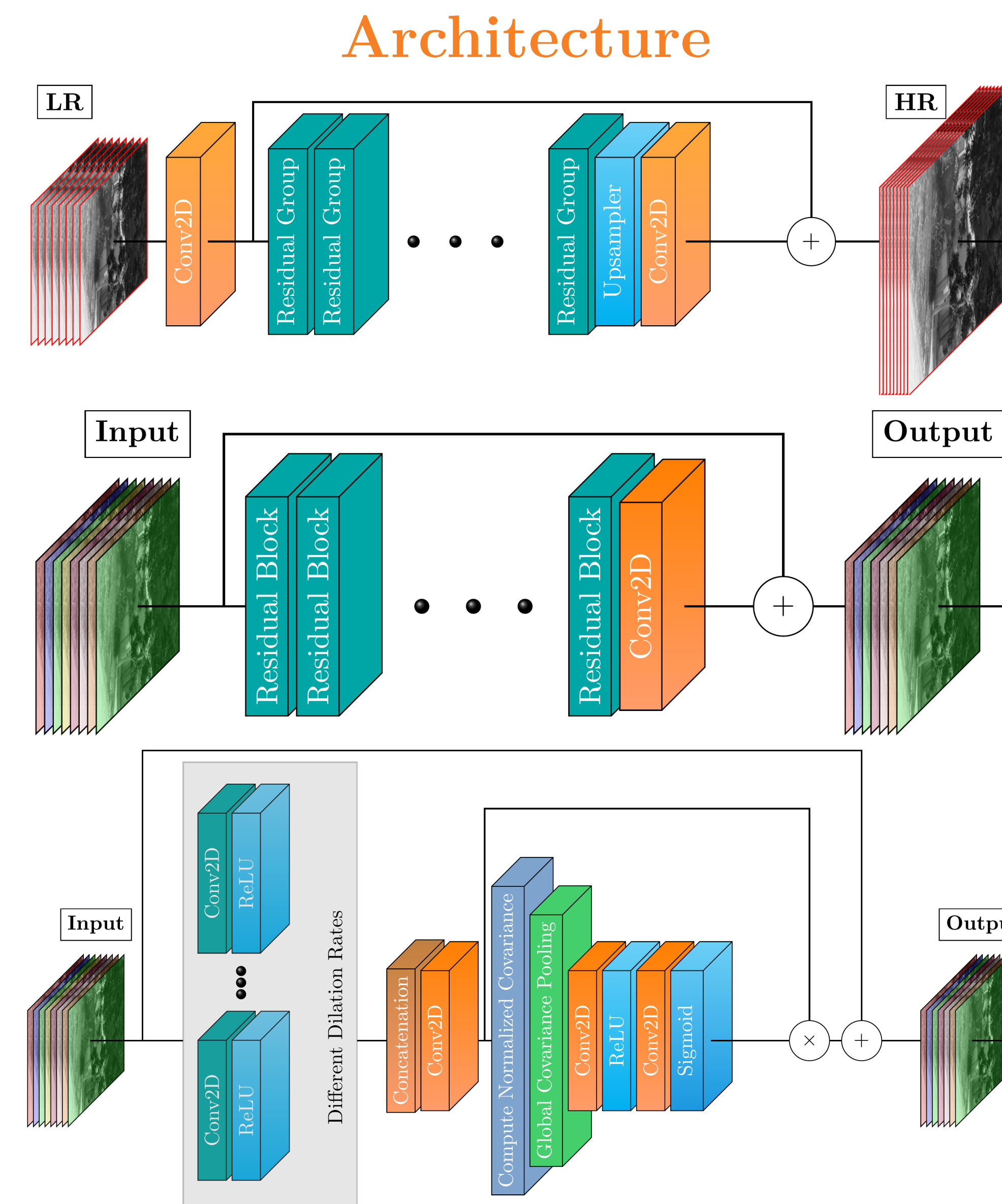


Figure 1: The residual body (top), residual group (middle), and residual block (bottom) of the AVRFN architecture.

This permits the channel attention to capture correlations higher than the first order.

Results

Model Variant	SOCA	Different Dilation Rates Module	Dilated Convolutions
DDRR		✓	✓
RRSOCA	✓		
CRCAN			✓
AVRFN	✓	✓	✓
RCAN			

Table 1: The model variants used during the ablation study.



Figure 2: Examples of (left column) downsampled images from (top row) low-resolution, (middle row) medium-resolution, and (bottom row) high-resolution thermal cameras, their $\times 4$ upscaled counterparts (middle column), and the ground truth (GT) (right column).

Test Set	Model	Scale	Parameters	PSNR	SSIM
AXIS Domo P1290	RRSOCA	4	1661377	25.487	0.691
	DDRR	4	2839873	25.458	0.691
	CRCAN	4	2839873	25.491	0.692
	RCAN	4	1661377	25.239	0.682
	AVRFN	4	1917313	25.368	0.685
AXIS Q2901	RRSOCA	4	1661377	28.167	0.802
	DDRR	4	2839873	28.159	0.801
	CRCAN	4	2839873	28.189	0.802
	RCAN	4	1661377	27.923	0.795
	AVRFN	4	1917313	27.990	0.797
FLIR FC-6320	RRSOCA	4	1661377	31.978	0.867
	DDRR	4	2839873	31.985	0.867
	CRCAN	4	2839873	32.002	0.867
	RCAN	4	1661377	31.756	0.861
	AVRFN	4	1917313	31.824	0.864
TDAT	RRSOCA	4	1661377	28.388	0.641
	DDRR	4	2839873	28.427	0.645
	CRCAN	4	2839873	28.426	0.645
	RCAN	4	1661377	28.271	0.636
	AVRFN	4	1917313	28.298	0.637
KAIST	RRSOCA	4	1661377	37.977	0.949
	DDRR	4	2839873	37.456	0.918
	CRCAN	4	2839873	37.573	0.922
	RCAN	4	1661377	37.089	0.938
	AVRFN	4	1917313	37.827	0.943

Table 2: Ablation study results showing peak signal-to-noise ratio (PSNR) and structural similarity index measure (SSIM) scores on the Thermal Image Super-Resolution (TISR), FLIR Thermal Dataset for Algorithm Training (TDAT), and KAIST datasets for several model variants.

We conducted an ablation study to evaluate how CTD, our different dilation rates module, and SOCA interact in an SR network. In Table 2, we show the evaluation results for model variants displayed in Table 1. The qualitative results of our best performing model variant (CRCAN) are shown in Fig. 2. Each of our model variants performs better than the strong RCAN baseline. Additionally, SOCA shows clear performance gains. However, CRCAN, which only has dilated convolutions, attains the greatest performance.

Acknowledgements

The authors acknowledge TACC at the University of Texas at Austin for providing software, computational, and storage resources that have contributed to the research results reported within this project.



UNIVERSITY OF
TEXAS
ARLINGTON